

Detailed Modeling and Evaluation of a Scalable Multi-level Checkpointing System

Kathryn Mohror, Adam Moody, Greg Bronevetsky, and Bronis R. de Supinski
Lawrence Livermore National Laboratory
{kathryn, moody20, bronevetsky, bronis}@llnl.gov



Abstract—High-performance computing (HPC) systems are growing more powerful by utilizing more components. As the system mean time before failure correspondingly drops, applications must checkpoint frequently to make progress. However, at scale, the cost of checkpointing becomes prohibitive. A solution to this problem is multilevel checkpointing, which employs multiple types of checkpoints in a single run. Lightweight checkpoints can handle the most common failure modes, while more expensive checkpoints can handle severe failures. We designed a multilevel checkpointing library, the Scalable Checkpoint/Restart (SCR) library, that writes lightweight checkpoints to node-local storage in addition to the parallel file system. We present a probabilistic Markov model of SCR's performance. Our model predicts that on future large-scale systems, SCR can lead to a gain in machine efficiency of up to 35%, and reduce the load on the parallel file system by a factor of two.

1 INTRODUCTION

Although supercomputing systems use high quality components, they become less reliable at larger scales because increased component counts increase overall fault rates. HPC applications can encounter mean times between failures (MTBFs) of hours or days due to hardware breakdowns [1] and soft errors [2]. For example, the 100,000 node BlueGene/L system at Lawrence Livermore National Laboratory (LLNL) experiences an L1 parity error every 8 hours [3] and a hard failure every 7-10 days. Exascale systems are projected to fail every 3-26 minutes [4], [5]. Most applications tolerate failures by periodically saving their state to reliable storage *checkpoint* files. Upon failure, an application can restart from a prior state by reading in a checkpoint.

This article has been authored by Lawrence Livermore National Security, LLC under Contract No. DE-AC52-07NA27344 with the U.S. Department of Energy. Accordingly, the United States Government retains and the publisher, by accepting the article for publication, acknowledges that the United States Government retains a non-exclusive, paid-up, irrevocable, world-wide license to publish or reproduce the published form of this article or allow others to do so, for United States Government purposes.

Checkpointing to a parallel file system is expensive at large scale. A single checkpoint can take tens of minutes [6], [7]. Further, large-scale computational capabilities have increased more quickly than I/O bandwidths. Typically, the limited bandwidth results from system design choices that optimize for system maintainability and availability.

Increasing failure rates due to increases in system scale require more frequent checkpoints. Increased system imbalance makes them more expensive. So, checkpointing is both more critical and less practical. Thus, large-scale applications require either more efficient checkpoint mechanisms or alternatives such as process replication, which have overheads over 100% [8].

Multilevel checkpointing [9], [10] uses multiple types of checkpoints that have different levels of resiliency and cost in a single application run to address this problem. The slowest but most resilient level writes to the parallel file system, which can withstand an entire system failure. Faster but less resilient levels use node-local storage, such as RAM, Flash or disk, and apply cross-node redundancy schemes. Most failures only disable one or two nodes, and multinode failures often disable nodes in a predictable pattern [11]. Thus, an application can usually recover from a less resilient checkpoint level, given carefully chosen redundancy schemes. Multilevel checkpointing allows applications to take frequent inexpensive checkpoints and less frequent, more resilient checkpoints, resulting in better efficiency and reduced load on the parallel file system.

We evaluate multilevel checkpointing in large-scale systems through a probabilistic Markov model. Our major contributions are:

- A Markov model of multilevel checkpointing;
- An exploration of modeled multilevel checkpointing performance on current and future systems;
- An evaluation of the viability of checkpointing to the parallel file system only upon job termination.

Overall, our results demonstrate that multilevel checkpointing significantly improves current methods. We show that it can increase system efficiency significantly, with gains up to 35% while reducing the load on the parallel file system by a factor of two.

The rest of this paper is organized as follows. Section 2 presents related work. In Section 3, we describe SCR, our Scalable Checkpoint/Restart library. Section 4 details our multilevel checkpoint model, Section 5 uses it to evaluate multilevel checkpointing on current and future systems. In Section 6, we extend our model to study the possibility of only writing checkpoints to the parallel file system when absolutely necessary.

2 RELATED WORK

Many models describe checkpoint systems [12], [13], [14], [15]. However, few have modeled multilevel checkpointing. Vaidya developed a Markov model for a two-level checkpoint system [16]. We extend Vaidya’s model to an arbitrary number of levels, each with its own checkpoint and recovery costs and failure rate. Our model also allows for sequential failures within a given computation interval.

Panda and Das extended Vaidya’s model to predict task completion probability. They assume a fixed number of spare resources and no repair [17]. In our model, we assume the system has an infinite pool of spare resources through repair of failed ones.

Gelenbe presented a Markov model for multilevel checkpointing [9]. He derived a formula for system efficiency from the Markov model steady state equations. However, he noted that an analytical solution for the optimum efficiency was intractable. Instead, we derive expressions for efficiency using a recursive method.

Researchers have combined checkpointing methods to lower overheads while maintaining resiliency [18], [19], [20]. However, to the best of our knowledge, SCR was the first implementation of multilevel checkpointing on large-scale, production systems.

3 THE SCR LIBRARY

SCR enables MPI applications to use node-local storage to attain high checkpoint and restart I/O bandwidth [11]. We derive its approach from two key observations. First, a job only needs its most recent checkpoint. As soon as it writes the next checkpoint, we can discard the previous one. Second, a typical failure only disables a small portion of the system.

Our SCR design leverages these observations by caching checkpoint files in storage local to the compute nodes instead of the parallel file system. SCR caches

only the most recent checkpoints, discarding an older checkpoint with each newly saved one. SCR can apply a redundancy scheme to the cache, so it can recover checkpoints after a failure disables some of the system. SCR periodically copies (flushes) a cached checkpoint to the parallel file system in order to protect against wider failures. However, a well-chosen redundancy scheme allows infrequent flushing of checkpoints.

LLNL uses SCR on Linux/x86-64/Infiniband clusters with RAM disks or solid-state drives (SSDs). The pF3D laser-plasma interaction code [21] has used SCR since late 2007 on the Hera, Atlas and Coastal systems at LLNL. We base our multilevel checkpoint model on SCR’s implementation and evaluate SCR’s performance using the model with parameters that represent SCR usage by pF3D.

4 MULTILEVEL CHECKPOINT MODEL

Our novel probabilistic model of multilevel checkpointing can predict the behavior of SCR given the factors that can affect its performance. This model can guide general use of multilevel checkpoint systems for current and future systems and motivate system designs that provide adequate overall reliability and efficiency. To model multilevel checkpointing systems, we make some simplifying assumptions that naturally introduce errors into the model’s predictions. However, these errors are relatively small. We now discuss our assumptions and their potential impact.

We assume that failures are independent. Thus, a failure within a job does not increase the probability of another failure within that job or future jobs. In reality, some failures are correlated. However, SCR is designed to mitigate effects of correlated failures. For example, it can avoid using failed nodes in a job allocation as those nodes may be likely to fail again.

We assume that checkpoints are taken at regular intervals throughout the job. While not always true, pF3D does checkpoint at regular intervals. We also assume costs to read and write checkpoints are constant throughout the job. However, read and write times actually vary, particularly when shared resources such as the parallel file system are used, which leads to some error in our model. We assume that the application recovers from the most recent viable checkpoint when a failure occurs. We do not model possible savings from using an older checkpoint that is also sufficient for recovery but available from faster storage. Thus, we may underestimate the possible performance of multilevel checkpointing.

We assume an infinite pool of spare nodes. When using SCR in practice, users often request extra nodes

in their job allocation; upon node failure SCR restarts the job using the extra nodes, ignoring those that failed. In general, the failure rate is less than the repair rate, so this assumption typically holds. In the absence of spare nodes, SCR copies the most recent checkpoint to the parallel file system and terminates the job; we model this capability in Section 6. Similarly, the model does not account for batch system allocation time limits. We assume a single level L checkpoint period completes within the allocation time limit. In practice, SCR handles batch limits by copying the most recent checkpoint to the parallel file system before the allocation expires. These assumptions cause our model to overestimate performance if SCR must copy checkpoints to the parallel file system.

4.1 Model Overview

In a multilevel checkpointing system, each of L checkpointing mechanisms is a *level*, for which level 1 checkpoints are the least expensive and resilient, while level L checkpoints are the most expensive and resilient. In our model, we assume that a checkpoint at level k can be used to recover from a superset of the failure modes that are recoverable using checkpoints at levels less than k . A *level k failure* refers to a failure severe enough that we require a checkpoint at level $i \geq k$ for recovery. A *level k recovery* restores an application using a checkpoint saved at level k . A multilevel checkpointing system alternates between different types of checkpoints. Since more severe failures happen less frequently, the system records zero or more level k checkpoints for every level $k + 1$ checkpoint.

In our Markov model of a multilevel checkpointing system, nodes represent application states and edges represent the transitions between states. We annotate each edge with the probability that the application will transition from the source state to the destination state and with cost information such as the time spent in the source state given that the transition is taken. Our model has *computation* and *recovery* states. Computation states represent periods of application computation followed by a checkpoint. Recovery states represent the process of restoring an application from a checkpoint saved previously.

Figure 1 presents our model's basic structure. The white states in the top row are computation states, and the single blue state at the bottom is a recovery state. We label each computation state by the checkpoint level with which it terminates and the recovery state by the checkpoint level that it uses to restore the application. If no failures occur during application

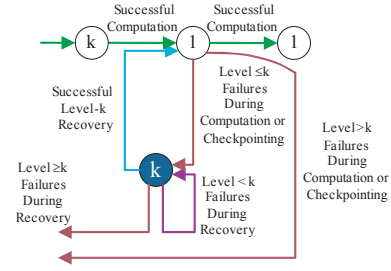


Fig. 1. Basic structure of multilevel Markov model

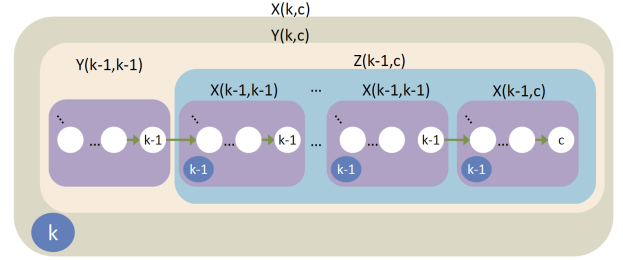


Fig. 2. Hierarchical structure of Markov model

execution or checkpointing, the application transitions from one computation state to the next. If a failure occurs, the application transitions to the recovery state corresponding to the most recent checkpoint capable of recovering from the failure. For example, if a failure at level i and $i \leq k$ while in the middle computation state in Figure 1, the system transitions to recovery state k , which restores the application using the checkpoint that was written at the end of the previous computation state. However, if $i \geq k$, the system must transition to a recovery state that corresponds to an older checkpoint saved at a higher level.

If no failures occur during recovery, the application transitions to the computation state that follows the checkpoint used for recovery. If a failure at level $i < k$ occurs while in a level k recovery, we assume the current recovery state must be restarted. However, if $i \geq k$, the application must transition to a higher-level recovery state. We assume a level L recovery can be restarted to recover after a failure at any level.

We exploit the recursive structure of our model to develop recurrence equations that we efficiently solve for the expected run time. As Figure 2 shows, we can build a full model by recursively composing three basic blocks, which we label $X(k,c)$, $Y(k,c)$, and $Z(k,c)$, where $k, c \in 1, 2, \dots, L$. An $X(k,c)$ block consists of a $Y(k,c)$ block and a base state for recovery at level k , R_k . A $Z(k,c)$ block consists of a series of $X(k,k)$ blocks and a terminating $X(k,c)$ block. When $k > 1$, a

Symbol	Definition
L	Number of checkpoint levels modeled
v_k	Number of level k checkpoints within each level $k + 1$ period
t	Length of compute interval before the application initiates a checkpoint
c_k	Time to record a level k checkpoint
r_k	Time to complete a level k recovery
λ_k	Average rate of level k failures assuming Poisson distributions

TABLE 1
Model parameters

$Y(k, c)$ block consists of either a single $Y(k-1, c)$ block or a $Y(k-1, k-1)$ block followed by a $Z(k-1, c)$ block. Finally, when $k = 1$, a $Y(k=1, c)$ block is a base state corresponding to a computation state that terminates with a checkpoint at level c . The parameter c is the checkpoint level taken by the last compute state in a block and k is the level of a block. An instance of $X(L, L)$ represents a level L interval.

We use the definitions listed in Table 3 to parameterize our multilevel checkpoint model. A technical report includes complete derivations of the results in Sections 9.3 to 9.6 [22]. In the interest of space, we only report a high level description here.

4.2 Base States

For the base computation and recovery states, p_0 is the probability that the application executes for some time, t_0 , without encountering a failure. For $k \in 1, 2, \dots, L$, the probability that the first failure during this period occurs at level k is p_k and t_k is the expected run time before encountering that failure. With T representing the time for spent in the state before exiting, and assuming an exponential distribution, the expressions for $p_0(T)$ and $t_0(T)$ evaluate to $p_0(T) = e^{-\lambda T}$ and $t_0(T) = T$, and for $k \in 1, 2, \dots, L$, $p_k(T)$ and $t_k(T)$ evaluate to

$$p_k(T) = \frac{\lambda_k}{\lambda} (1 - e^{-\lambda T}), \quad t_k(T) = \frac{1 - (\lambda T + 1) \cdot e^{-\lambda T}}{\lambda \cdot (1 - e^{-\lambda T})},$$

where $\lambda = \lambda_1 + \lambda_2 + \dots + \lambda_L$.

A $Y(k=1, c)$ block is a base computation state in which the application executes for an interval of length t and then writes a checkpoint at level c , which requires a time of c_c . From the formulas above, $p_{Y0} = p_0(t + c_c)$ and $t_{Y0} = t_0(t + c_c)$, and for $i \in 1, 2, \dots, L$, $p_{Yi} = p_i(t + c_c)$, and when $p_{Yi} > 0$, $t_{Yi} = t_i(t + c_c)$.

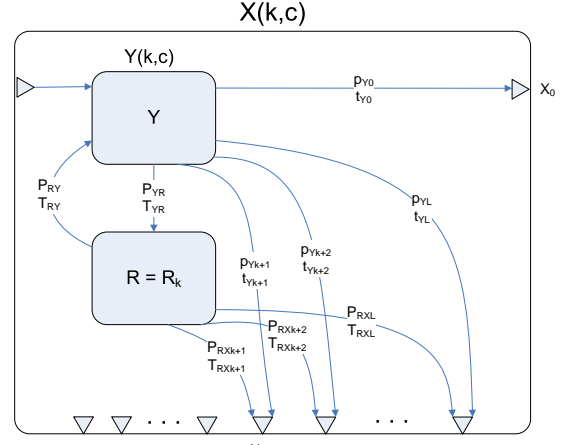


Fig. 3. Simplified diagram of $X(k, c)$

While in a recovery base state at level k , the system is recovering from a failure using a checkpoint saved at level k , which requires a time of r_k . We find that the probability of exiting with no failures is $p_{R0} = p_0(r_k)$ and the time to exit with no failures is $t_{R0} = t_0(r_k)$. For $i \in 1, 2, \dots, L$, the probability of exiting on a failure at level i is $p_{Ri} = p_i(r_k)$, and when $p_{Ri} > 0$, the time before exiting on failure at level i is $t_{Ri} = t_i(r_k)$.

4.3 The $X(k, c)$ block

An $X = X(k, c)$ block internally consists of a $Y = Y(k, c)$ block and a recovery state at level k , $R = R_k$. To simplify the final expressions, we merge groups of related transitions into single transitions. We show the merged transitions in Figure 13(b).

Y transitions to the recovery state R for any failure scenario that requires a recovery level at k or less. We merge each of these transitions into a single transition that has probability of P_{YR} and an expected run time of T_{YR} . Once in R , a transition away from R eventually happens, provided that $\sum_{i=0}^k p_{Ri} < 1$. However, one or more loops back to R may occur before transitioning away. We merge the transitions from R to Y ; its probability is P_{RY} and its expected run time is T_{RY} as shown in Figure 13(b).

Failures at levels $i \geq k$ cause a transition out of R_k to a higher level recovery state. The transitions from R have probability P_{RXi} and expected run time T_{RXi} . However, if $k = L$, then recovery is restarted upon failure at any level so for each $i \in 1, 2, \dots, L$, $P_{RXi} = 0$. While in a recovery state at level $k < L$, the system

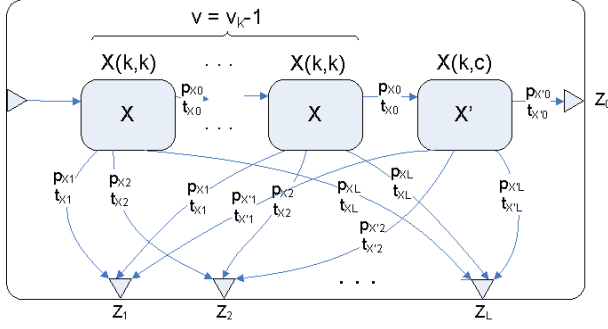


Fig. 4. The $Y(k, c)$ State

transitions to a recovery state at level $k+1$ if a level k or level $k+1$ failure occurs. Otherwise, for the occurrence of a failure at level i , where $i > k+1$, a transition is made to a recovery state at level i .

4.4 The $Z(k, c)$ block

A $Z = Z(k, c)$ block only exists when $v_k > 0$. It consists of a chain of $X = X(k, k)$ blocks of length $v_k - 1$ followed by a $X' = X(k, c)$ block, as Figure 15 shows. We define $v = v_k - 1$.

The probability of successfully transitioning from the Z block to the first computation state of the next block is the probability that v consecutive successful transitions from X blocks are followed by one successful transition from the X' block, $p_{Z0} = (p_{X0})^v \cdot p_{X'0}$. When $p_{Z0} > 0$, the expected time to make this transition is $t_{Z0} = v \cdot t_{X0} + t_{X'0}$. The total probability to leave Z for a recovery state at level i is the sum of the probabilities corresponding to each of the possible paths from the substates of Z .

4.5 The $Y(k, c)$ block

A $Y(k, c)$ block is built using three different constructions depending on the values of k and (when $k > 1$) v_{k-1} . If $k = 1$, then $Y(k, c) = Y(k = 1, c)$, which is a base computation state. The probability and expected run time vectors for this state can be directly computed as described in Section 9.3.

If $k > 1$ and $v_{k-1} = 0$, then $Y = Y(k, c)$ consists of a single $Y' = Y(k - 1, c)$ block. We compute the probabilities and expected run times to transition from Y given the probabilities and expected run times to transition from Y' as $p_{Y0} = p_{Y'0}$ and $t_{Y0} = t_{Y'0}$, and, for each level $i \in 1, 2, \dots, L$, $p_{Yi} = p_{Y'i}$ and $t_{Yi} = t_{Y'i}$.

If $k > 1$ and $v_{k-1} > 0$, then $Y = Y(k, c)$ consists of a starting $Y' = Y(k - 1, k - 1)$ block followed by a $Z = Z(k - 1, c)$ block, as Figure 16 shows.

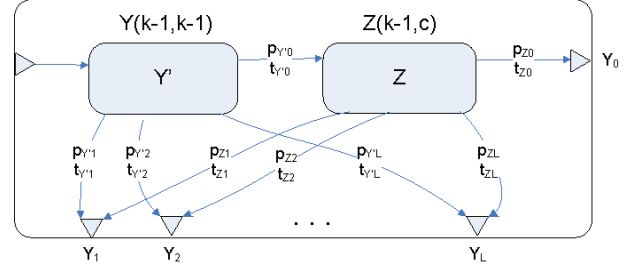


Fig. 5. The $Y(k, c)$ State for $k > 1$ and $v_{k-1} > 0$

The probability that a successful transition from Y occurs is the probability that both Y' and Z transition successfully, $p_{Y0} = p_{Y'0} \cdot p_{Z0}$ and the expected time for this transition is $t_{Y0} = t_{Y'0} + t_{Z0}$. The total probability to leave Y for a recovery state at level i is the sum of the probabilities for each path.

4.6 Model Metrics

We consider two key metrics: *efficiency* and *parallel file system load*. We define *efficiency* as the ratio of *idealTime* to *expectedTime*, where *idealTime* is the minimum run time assuming the application spends no time checkpointing and encounters no failures, while *expectedTime* is the expected run time that the model predicts for a set of parameters. This metric indicates how much time is lost to checkpointing, including recovery from failures.

To compute efficiency, we parameterize the model with a set of checkpoint levels including their checkpoint and recovery costs, failure rates, and time between checkpoints. We then compute the expected time to complete a level L period. The value of t_{X0} for the $X(L, L)$ state is the *expectedTime* to complete a level L period. The *idealTime* is the total number of compute intervals multiplied by the length of each interval.

To judge the impact on the parallel file system for a particular model configuration, we consider the expected time between writing consecutive checkpoints to the parallel file system. We define the *load* on the parallel file system to be the inverse of *expectedTime*.

5 MODEL EXPLORATION

We used our model to explore the behavior of SCR under varying conditions. We show the benefits of multilevel over single-level checkpointing as failure rates and parallel file system characteristics change.

We show predictions of pF3D efficiency to those observed in real runs on Coastal and Atlas in Table 5. The data show that the model's predictions are within

System	Expected Efficiency	Observed Efficiency	Duration of Observation
Coastal	95.2%	94.68%	716,613 node-hours
Atlas	96.7%	92.39%	553,829 node-hours

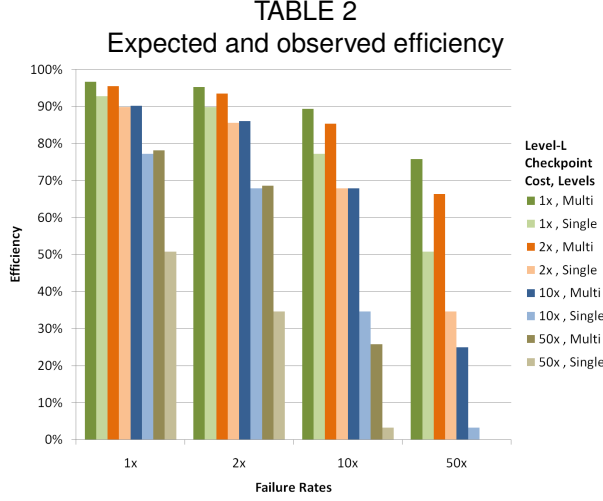


Fig. 6. Optimal efficiency for single- and multilevel checkpointing

a few percent of our observations. Despite limitations due to the significant time required to gather data, these results demonstrate that our model is accurate.

We now use the model to explore multilevel checkpointing in a more general context. For checkpoint costs, we use the times observed for checkpointing pF3D on Coastal using three different checkpointing levels, with costs of 0.5 seconds, 4.5 seconds, and 1052 seconds, with recovery costs equal to checkpoint costs. Using collected failure data for pF3D on Coastal, we use rates in failures per job-second of $2 \cdot 10^{-7}$ for level 1, $1.8 \cdot 10^{-6}$ for level 2, and $4 \cdot 10^{-7}$ for level 3.

As future systems become larger, failure rates are expected to increase, and as the system memory size grows faster than the performance of the parallel file system, the cost of accessing the parallel file system is expected to increase. To explore these effects, we increase the base failure rates and the level L checkpoint costs by factors of 2, 10, and 50. We do not adjust the costs of lower-level checkpoints, since the performance of node-local storage is expected to scale with system size. For each combination, we identified the compute interval and the level 1 and level 2 checkpoint counts that provide the highest efficiency. We compare this to single-level checkpointing, for systems with only a parallel file system available.

Figure 6 presents the efficiency achieved for each configuration, and Figure 7 shows the reduction in load

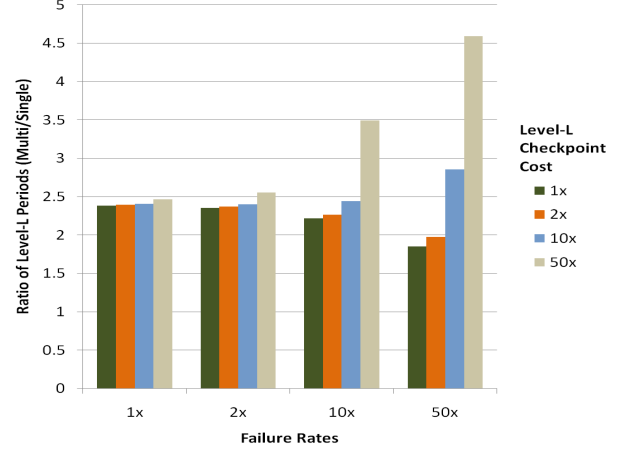


Fig. 7. Parallel file system load reduction

on the parallel file system. We label the results for the multilevel system as “Multi” and those for the single-level system as “Single.” The groupings of bars along the x-axis correspond to failure rates that are 1, 2, 10, or 50 times the base values. Within each grouping, we increase the cost of the level L checkpoint by 1, 2, 10, and 50 times the base value.

In all cases, the multilevel system results in higher efficiencies and increases the time between checkpoints to the parallel file system. Moreover, both advantages increase with either increasing failure rates or higher parallel file system costs. The gain in machine efficiency ranges from a few percent up to 35%, and, as Figure 7 shows, the load on the parallel file system is reduced by a factor ranging from 2x-4x. Thus, compared to single-level checkpointing, multilevel checkpointing simultaneously increases efficiency while reducing load on the parallel file system. These results highlight the benefits of multilevel checkpointing on current and future systems.

Overall, we find that multilevel checkpointing is essential for future systems. Even with systems that are 50x less reliable, a three level checkpointing system achieves efficiencies over 75%, as long as we maintain relative parallel file system performance. On the other hand, we find that we cannot tolerate higher failure rates if the cost to access the parallel file system also increases. In particular, if systems become 50x less reliable and the cost of saving application state to the parallel file system rises by 10x, a three level checkpointing system only achieves 26% efficiency.

6 SCAVENGING OF LAST CHECKPOINT

We now extend our model to the scenario in which we only write level L checkpoints upon job termination, and not periodically during the run. This scenario could be advantageous on systems that have high costs for level L checkpoints. Also some systems do not support application restart on its current allocation when a failure occurs; instead, the job must be restarted in a new allocation. In this case, upon application termination, a multilevel checkpoint system should write the last complete checkpoint to the parallel file system. This approach only incurs the high overhead of level L checkpoints when strictly necessary. However, this carries the risk of the level L checkpoint failing because the lower-level checkpoints were corrupted.

We refer to the process of pushing a checkpoint set to the parallel file system upon allocation termination as *scavenging a checkpoint*. Scavenging only occurs on application failure. Level $k < L$ checkpoints are taken at regular intervals during the application run and cached on the compute nodes.

We make several assumptions in our model. In the event of any failure, the application does not restart; instead, the scavenge process begins. On a level k failure, we attempt to scavenge the most recent checkpoint at level $i \geq k$. If a failure occurs at level $i \geq k$ while scavenging, we attempt to scavenge the most recent checkpoint at level j or greater, where $j > i$. If a level L failure occurs while scavenging, the scavenge operation fails. When the application is restarted in a new allocation, it must roll back to the last level L checkpoint taken before the prior allocation. Finally, we estimate the time to scavenge the level k checkpoint as the time for a level L checkpoint. In reality, we may incur a small additional overhead to rebuild the level k checkpoint depending upon the failure mode.

We evaluate the effectiveness of scavenging by computing the expected efficiency of the job, as defined in Section 4.6. However, we compute it differently, because the job will terminate on any failure and not attempt restarts. We define *efficiency* as the ratio of the expected amount of work done by the application before its expected termination time either on success, w_{X0} , on scavenge, w_{XS} , or on failure without a successful scavenge, w_{Xi} . In our model, the work done by the application is the checkpoint interval, t , multiplied by the number of intervals completed successfully; work excludes the time for any checkpointing activities. The expected termination time is the time to exit on success, t_{X0} , on scavenge, t_{XS} , or on failure without successful scavenge, t_{Xi} , and includes the work time, level $k \leq L$

checkpointing time, and scavenge time.

$$\text{efficiency} = \frac{p_{X0} \cdot w_{X0} + p_{XS} \cdot w_{XS} + p_{Xi} \cdot w_{Xi}}{p_{X0} \cdot t_{X0} + p_{XS} \cdot t_{XS} + p_{Xi} \cdot t_{Xi}}$$

We base our revised efficiency computation on the amount of work performed in each state. On failure, if scavenging succeeds, we compute the expected amount of work completed before the failure. If scavenging fails, we compute the expected amount of work completed in failure transitions only. A transition from state on success means that no failures occurred and the amount of work completed is the work accumulated in any substates of the current state.

If a level k failure occurs and the current state contains a level $i \geq k$ checkpoint, we can scavenge a checkpoint. We call this *transitioning on scavenge*. In this case, the amount of work completed is the amount of work accumulated in substates before the failure. A transition on failure occurs if the current state does not contain a level $i \geq k$ checkpoint or a level $i \geq k$ failure occurs while scavenging. On failure transitions at level $i > k$, the amount of work completed in the current state is zero. In the extreme case, if no completed checkpoints can be scavenged and the application must roll back to the last level L checkpoint, the total expected amount of work completed is zero.

We make several changes to our model. Recovery states within $X(k, c)$ are now scavenge states. Transitions from each compute state at failure levels $i \leq k$ are *scavenge transitions*. Instead of maintaining the probability, expected total time, and expected amount of work at each failure level for scavenge transitions, we compute the overall expected values. Unless otherwise specified, the probabilities and times for exiting states on failure and success are the same as in Sections 9.3 to 9.6.

6.1 Y States

The simplest case is a $Y(k = 1, c)$ block, which contains a single compute interval of duration t followed by a level c checkpoint. Upon successful transition from this state, the expected amount of work completed is t . Upon transition from Y on failure for any level $i \in 1, 2, \dots, L$, $w_{Yi} = 0$ because no work interval and checkpoint completed successfully. We cannot transition from Y on scavenge, because a $Y(k = 1, c)$ state does not contain a scavenge state. Thus, on transitioning from Y , $p_S = 0$, $t_S = 0$, and $w_S = 0$.

If $v_{k-1} = 0$ and $k > 1$, Y consists of a single enclosed $Y(k-1, c)$ state, Y' . The work completed for a transition from Y on success or failure is $w_{Y0} = w_{Y'0}$

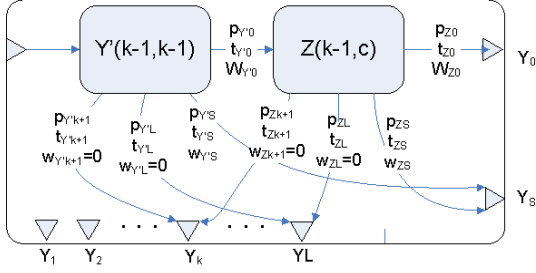


Fig. 8. The $Y(k, c)$ State for $v_{k-1} > 0$

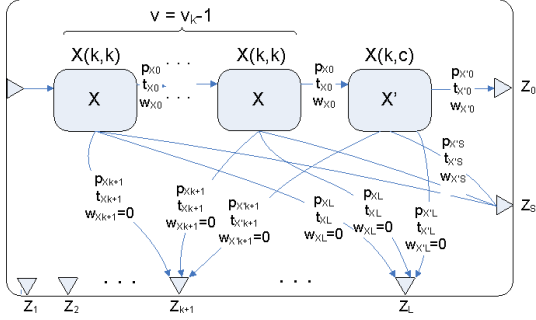


Fig. 9. $Z(k, c)$ State

and $w_{Y_i} = w_{Y'_i} = 0$. For a scavenge transition, we have $p_{YS} = p_{Y'S}$, $t_{YS} = t_{Y'S}$, and $w_{YS} = w_{Y'S}$.

If $v_{k-1} > 0$ and $k > 1$, Y consists of a $Y' = Y(k-1, k-1)$ state and a $Z = Z(k-1, c)$ state (Figure 17(c)). On successful transition from Y , w_{Y0} is the sum of the work done in Y' and Z , $w_{Y0} = w_{Y'0} + w_{Z0}$. Upon a transition from Y on failure at levels $i \in 1 \dots L$, $w_{Y_i} = 0$. The probability of transitioning on scavenge p_{YS} is the probability that the job exits on scavenge in Y' or completes Y' successfully and exits on scavenge from Z . Similarly, the time and work completed, t_{YS} and w_{YS} depend on the probabilities of transitioning from Y' and Z on success or scavenge.

6.2 Z State

A $Z = Z(k, c)$ state consists of $v = v_k - 1$ consecutive $X = X(k, k)$ states, followed by a $X' = X(k, c)$ state (Figure 18). The expected amount of work that is performed in Z on successful transition from the state is $w_{Z0} = w_{X0} \cdot v + w_{X'0}$. Upon transition from Z on failure for $i \in 1 \dots L$, the amount of work completed is $w_{X_i} = 0$.

6.3 X State

An $X = X(k, c)$ state consists of a $Y = Y(k, c)$ state and a scavenge state σ_k . A failure at level $i \leq k$ in Y

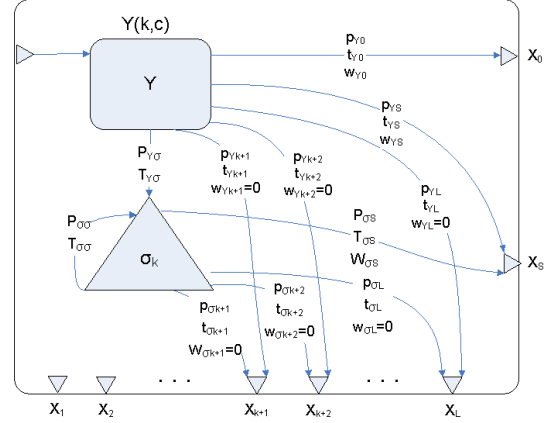


Fig. 10. $X(k, c)$ State

causes a transition from Y to σ_k . For any failures at levels $i < k$ that occur in σ_k , the scavenge operation restarts. Failures at levels $i \geq k$ in Y or $i \geq k$ in σ_k , cause a transition from X on failure.

The probability and time for a successful transition from X are $p_{X0} = p_{Y0}$ and $t_{X0} = t_{Y0}$. The expected amount of work that is performed in X on successful transition from the state is the amount of work that was performed in Y , $w_{X0} = w_{Y0}$.

A failure at level $i \leq k$ in Y causes a transition to the scavenge state σ_k . The amount of work done on transition from Y to σ_k is $w_{Y\sigma_k} = 0$. As stated previously, we estimate the time for a successful scavenge operation as the time for a level L checkpoint, c_L . Failures at levels $i < k$ restart the scavenge operation. The probability and time before a restart of the scavenge operation are $P_{\sigma\sigma}$ and $T_{\sigma\sigma}$. No work is done in σ_k , so $w_{\sigma_i} = 0$.

The probability and time for leaving σ_k on success or failure for $i \in 0, 1, \dots, L$ depend on the probabilities for exiting σ_k successfully after a number of restarts. The amount of work accomplished is $W_{\sigma_i} = w_{\sigma0}$. However, if $P_{\sigma\sigma} = 1$ then $P_{\sigma i}$, $T_{\sigma i}$, and $W_{\sigma i}$ are zero.

When transitioning from X on failure at level $i = k$, the probability, time, and work only depend on the scavenge state. However, when transitioning from X on failure at level $i > k$, the values depend on the Y and scavenge states.

6.4 Model Predictions of Benefits of Scavenge

We now use our model to predict the performance of scavenging. In our experiments, we use the overhead and failure rates for the Coastal cluster at LLNL, given in Section 5. We first explore the relationship between application efficiency and the compute interval. We

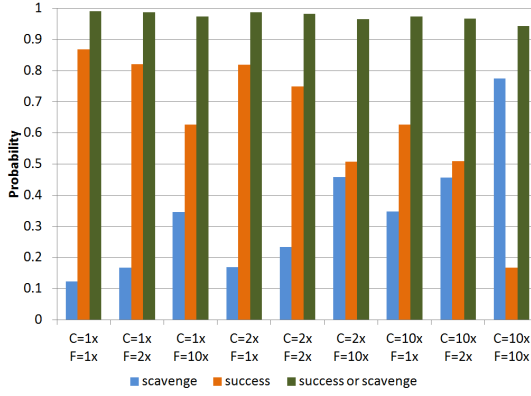


Fig. 11. Application success or scavenge probabilities

vary the failure rates and overhead of writing to the parallel file system by $1\times$, $2\times$, and $10\times$ as described previously in Section 5. The trend for today's failure rates is a broad range of compute intervals with near-optimal efficiencies (89%, 95%, and 96%) that level off near 80%, 90%, and 92% as overheads increase by $1\times$, $2\times$, and $10\times$ and checkpoint intervals increase up to 10,000 seconds. At failure rates that are $2\times$ greater than today's, the ranges of near-optimal efficiencies are shorter, peak at 83%, 93%, and 95%, and level off at 74%, 83%, and 87% efficiencies for increasing checkpoint overhead and checkpoint interval up to 10,000 seconds. However, at failure rates of $10\times$, the curves reach lower optimal efficiencies (55%, 81%, and 88%) and drop off more sharply with increasing compute interval to 46%, 67, and 70% efficiencies.

We predict the probability that an application will exit with its last checkpoint written to the parallel file system, either with no failures or on scavenge. For this experiment, we fix the probability of exiting the scavenging operation at 80%, based on observed successful scavenge rates with pF3D. We show the results for varying parallel file system overhead and failure rates in Figure 11. The probability of exiting with no failures starts at 86% and decreases with both increasing parallel file system and failure rates to a low of 17% when the overhead and failure rates are increased to $10\times$. In contrast, the probabilities of exiting on scavenge increase with increasing overhead and failure rates, from 12% at today's overhead and rates to 77% when they are increased $10\times$. The combined probability of exiting with no failures or on scavenge is 99% on today's systems and falls slightly with increasing overhead and failure rates to a low of 94% when the overhead and rates are at $10\times$. The remaining 1-6% are the probabilities that the application exits

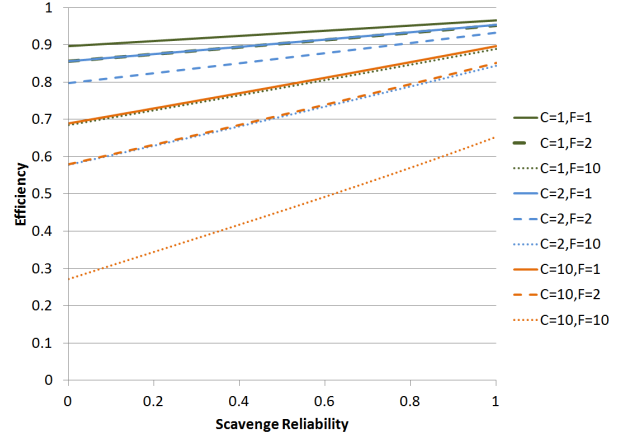


Fig. 12. Reliability of scavenge operation and efficiency

without being able to transfer the final checkpoint to the parallel file system.

We use our model to predict the reduction in load on the parallel file system when using scavenging compared to single-level checkpointing. Here, the load is the expected time between writes to the parallel file system. Scavenging greatly reduces the load on today's systems, by $20\times$. As parallel file system overhead increases, the benefit of scavenging decreases; however the lowest benefit in our experiments is still $10\times$ reduction in load. In general, the benefit increases with increasing failure rate, reaching a maximum reduction of $60\times$ for failure rates at ten times today's values.

In Figure 12, we explore the relationship between the reliability of the scavenge operation and application efficiency. We fix the probability of completing the scavenge operation to a given value instead of computing the probability as given in Section 10.3. Generally, job efficiencies are high for systems with today's failure rates, because the probabilities of exiting the job successfully without needing to scavenge are relatively high. However, with increasing failure rates and parallel file system overheads, the reliability of scavenging affects job efficiency more drastically. We find that on future systems with higher failure rates, a highly reliable scavenge operation can increase job efficiency by as much as $2.4\times$.

Overall, our results show that a scavenge mechanism can dramatically extend the range of systems for which checkpoint/restart remains a viable resilience strategy.

7 CONCLUSIONS

We presented models of multilevel checkpointing, which we validated against results with the Scal-

able Checkpoint/Restart (SCR) library. SCR combines checkpointing to stable storage with lower-overhead, less-resilient checkpoint types, e.g., copying checkpoints to memory on other nodes. Our novel, hierarchical Markov models predict the performance of multilevel checkpointing systems based on system reliability and checkpoint cost. This model can guide users in selecting the best checkpointing parameters for their application. Our analysis with this model demonstrates that multilevel checkpointing significantly improves system efficiency, particularly as failure rates and relative parallel file system checkpoint costs increase. We find that we can still achieve 85% efficiency even if systems become $50\times$ less reliable. Further, multilevel checkpointing simultaneously reduces the load on the parallel file system by more than a factor of two.

We explored the impact of checkpoint scavenging, a key extension to multilevel checkpointing. We found that scavenging results in high efficiencies even when overheads of the parallel file system are increased by $2\times$ and $10\times$. Further, the extremely high probability (99%) that the final checkpoint of the job can be transferred to the parallel file system means that work completed in the failed job is not lost and the job can be restarted in a new allocation. Scavenging also has the benefit of only writing to the parallel file system when absolutely necessary. Our model predicts that scavenging can reduce the load on the parallel file system by as much as $20\times$ on today's systems, and up to $60\times$ on future systems. We also found that a highly reliable scavenge operation can increase job efficiencies by up to $2.4\times$ on systems with higher failure rates and parallel file system overheads.

REFERENCES

- [1] B. Schroeder and G. A. Gibson, "A Large-Scale Study of Failures in High-Performance Computing Systems," in *Proceedings of the International Conference on Dependable Systems and Networks (DSN)*, June 2006, pp. 249–258.
- [2] S. E. Michalak, K. W. Harris, N. W. Hengartner, B. E. Takala, and S. A. Wender, "Predicting the Number of Fatal Soft Errors in Los Alamos National Laboratory's ASC Q Supercomputer," *IEEE Transactions on Device and Materials Reliability*, vol. 5, no. 3, pp. 329–335, September 2005.
- [3] J. N. Glosli, K. J. Caspersen, J. A. Gunnels, D. F. Richards, R. E. Rudd, and F. H. Streitz, "Extending Stability Beyond CPU Millennium: A Micron-Scale Atomistic Simulation of Kelvin-Helmholtz Instability," in *Proceedings of the 2007 ACM/IEEE Conference on Supercomputing (SC)*, 2007, pp. 1–11.
- [4] B. Schroeder and G. Gibson, "Understanding Failure in Petascale Computers," *Journal of Physics Conference Series: SciDAC*, vol. 78, p. 012022, June 2007.
- [5] E. Vivek Sarkar, Ed., *ExaScale Software Study: Software Challenges in Exascale Systems*, 2009.
- [6] K. Iskra, J. W. Romein, K. Yoshii, and P. Beckman, "ZOID: I/O-Forwarding Infrastructure for Petascale Architectures," in *PPoPP '08: Proceedings of the 13th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, 2008, pp. 153–162.
- [7] R. Ross, J. Moreira, K. Cupps, and W. Pfeiffer, "Parallel I/O on the IBM Blue Gene/L System," Blue Gene/L Consortium Quarterly Newsletter, Tech. Rep., First Quarter, 2006.
- [8] R. E. Lyons and W. Vanderkulk, "The Use of Triple-Modular Redundancy to Improve Computer Reliability," *IBM Journal of Research and Development*, vol. 6, no. 2, pp. 200–209, 1962.
- [9] E. Gelenbe, "A Model of Roll-back Recovery with Multiple Checkpoints," in *Proceedings of the 2nd International Conference on Software Engineering (ICSE '76)*, 1976, pp. 251–255.
- [10] N. H. Vaidya, "A Case for Multi-Level Distributed Recovery Schemes," Texas A&M University, Tech. Rep. 94-043, May 1994.
- [11] A. Moody, G. Bronevetsky, K. Mohror, and B. R. de Supinski, "Design, Modeling, and Evaluation of a Scalable Multi-level Checkpointing System," in *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, SC'10*, November 2010, pp. 1–11.
- [12] J. W. Young, "A First Order Approximation to the Optimum Checkpoint Interval," *Communications of the ACM*, vol. 17, no. 9, pp. 530–531, 1974.
- [13] A. Duda, "The Effects of Checkpointing on Program Execution Time," *Information Processing Letters*, vol. 16, no. 5, pp. 221–229, 1983.
- [14] J. S. Plank and M. G. Thomason, "Processor Allocation and Checkpoint Interval Selection in Cluster Computing Systems," *Journal of Parallel Distributed Computing*, vol. 61, no. 11, pp. 1570–1590, 2001.
- [15] J. Daly, "A Higher Order Estimate of the Optimum Checkpoint Interval for Restart Dumps," *Future Generation Computer Systems*, vol. 22, no. 3, pp. 303–312, 2006. [Online]. Available: <http://www.sciencedirect.com/science/article/B6V06-4F490KH-6/2/6ebfa65591e5d0eb09e2ae5ae3b2ed44>
- [16] N. H. Vaidya, "A Case for Two-Level Distributed Recovery Schemes," in *Proceedings of the 1995 ACM SIGMETRICS Joint International Conference on Measurement and Modeling of Computer Systems (SIGMETRICS '95)*, 1995, pp. 64–73.
- [17] B. S. Panda and S. K. Das, "Performance Evaluation of a Two Level Error Recovery Scheme for Distributed Systems," in *4th International Workshop on Distributed Computing, Mobile and Wireless Computing (IWDC)*, 2002, pp. 88–97.
- [18] L. Silva and J. Silva, "Using Two-Level Stable Storage for Efficient Checkpointing," *IEEE Proceedings - Software*, vol. 145, no. 6, pp. 198–202, Dec 1998.
- [19] J. S. Plank and K. Li, "Faster Checkpointing with N+1 Parity," in *Twenty-Fourth International Symposium on Fault-Tolerant Computing (FTCS), Digest of Papers*, Jun 1994, pp. 288–297.
- [20] L. A. Bautista-Gomez, S. Tsuboi, D. Komatitsch, F. Cappello, N. Maruyama, and S. Matsuoka, "FTI: High Performance Fault Tolerance Interface for Hybrid Systems," in *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, SC'11*, 2011.
- [21] R. L. Berger, C. H. Still, E. A. Williams, and A. B. Langdon, "On the Dominant and Subdominant Behavior of Stimulated Raman and Brillouin Scattering Driven by Nonuniform Laser Beams," *Physics of Plasmas*, vol. 5, p. 4337, 1998.
- [22] A. Moody, G. Bronevetsky, K. Mohror, and B. R. de Supinski, "Detailed Modeling, Design, and Evaluation of a Scalable Multi-level Checkpointing System," Lawrence Livermore National Laboratory Technical Report, LLNL-TR-440491, Tech. Rep., 2010.

8 SUPPLEMENTAL MATERIAL

In this document, we give derivation details for our model in the original paper. Please refer to the paper for a high level description of the model.

9 MODEL DERIVATIONS

9.1 Probability and Expected Run Time Vectors

Upon transitioning into the first computation state of an $X(k, c)$, $Y(k, c)$, or $Z(k, c)$ block, the system will eventually either transition out of the block to the first computation state of the next block, or it will transition to an external recovery state at some level $k \in 1, 2, \dots, L$. We represent the probabilities and expected run times for each of these transitions in vectors \vec{p} and \vec{t} , which we compute for each block. Each vector has $L+1$ elements where the i -th element, for $i \in 0, 1, \dots, L$, is labeled p_i and t_i , respectively. The elements of \vec{p} are labeled p_i for $i \in 0, 1, \dots, L$, and the elements of \vec{t} are labeled t_i for $i \in 0, 1, \dots, L$. The variable p_0 represents the probability that a transition is made from a block to the first computation state of the next block, and t_0 represents the expected run time spent within the block given such a transition. Variables p_k , for $k \in 1, 2, \dots, L$ represent the probability that a transition is made from a block to an external recovery state due to a failure scenario requiring a recovery at level k , and element t_k represents the expected run time spent in the block given such a transition.

We define \vec{p} and \vec{t} similarly for the base recovery states. In this case, p_0 and t_0 represent the probability and expected run time of completing the recovery process without failure and transitioning to the computation state that follows the checkpoint used for recovery. Element p_k represents the probability of encountering a level k failure while in recovery, and t_k represents the expected run time before encountering such a failure.

Given the constituent components that a block is built from, along with the \vec{p} and \vec{t} vectors for each of those components, we can compute the \vec{p} and \vec{t} vectors for the block. Starting from the base computation and recovery states, we can compute the probability and expected run time vectors for all blocks in a given model structure.

We use the definitions listed in Table 3 to parameterize our multi-level checkpoint model. The number of checkpoint levels is represented by L . The number of level k checkpoints taken within each level $k+1$ period is represented by v_k , for $k \in 1, 2, \dots, L-1$. We assume the application checkpoints (to some level) after completing regular intervals of computation. We

represent the length of this compute interval by t . The parameter c_k represents the time required to write a level k checkpoint, and r_k represents the time required to restore an application using a level k checkpoint. In this work, we assume failures at each level follow a Poisson distribution, where λ_k represents the average failure rate at level k .

The multi-level checkpointing model relies on a number of definitions and notation details. These are listed in Table 3. Symbols t , c_k and r_k determine the time spent performing the basic computation and reliability operations, with t being the duration of a level 1 execution period, and c_k and r_k being the times to record a level k checkpoint and perform a level k recovery, respectively. L and V determine the structure of the multi-level checkpointing system, with L being the number of levels modeled and v_k being the number of level k periods inside each level $k+1$ period. Functions $f_i(t)$ and $F_i(T)$ determine the reliability of the modeled system, with $f_k(t)$ being the probability of suffering a level k failure at time t and $F_k(T)$ the probability of a level k failure in the time period $[0..T]$. Finally, the function $p_i(T)$ and $t_i(T)$ are the probability of and expected time until an application activity that takes T time is aborted due to a level i failure. $p_0(T)$ and $t_0(T)$ are the probability of and expected time until a successful completion of a T -duration activity, because no failures are encountered. These determine the probability of and expected time to transition from a given Markov state either successfully or unsuccessfully.

9.2 Basic Derivations

We will utilize the following formulas in our analysis. The following two formulas compute the sums of geometric series for when $x \neq 1$:

$$\sum_{i=0}^N x^i = \frac{1 - x^{(N+1)}}{1 - x}, \quad (1)$$

$$\sum_{i=1}^N i \cdot x^i = \frac{x - (N+1) \cdot x^{(N+1)} + N \cdot x^{(N+2)}}{(1-x)^2}. \quad (2)$$

When a block has multiple edges that transition to the same destination block, we apply the following formulas to merge those edges into a single logical edge. Assume a block has N edges with probabilities p_i and expected costs t_i for $i \in 1, 2, \dots, N$ that all transition to the same destination block. These edges can be combined into a single logical edge having

Symbol	Definition
L	Number of checkpoint levels being modeled
v_k	Number of level k checkpoints within each level $k + 1$ period
t	Length of compute interval before application initiates a checkpoint
c_k	Time to record a level k checkpoint
r_k	Time to complete a level k recovery
λ_k	Average rate of level k failures assuming Poisson distributions
$p_0(T)$	The probability of completing an activity that takes T time without experiencing any failures
$t_0(T)$	Given no failures occurred during an activity of duration T , the expected time to complete this activity.
$p_i(T)$	The probability of aborting an activity (computation, checkpointing or recovery) that takes T time due to level i failure
$t_i(T)$	Given that an activity of duration T was aborted by a level i failure, the expected time until this happens.

TABLE 3
Model parameters

probability P and expected cost T computed as

$$P = \sum_{i=1}^N p_i \quad (3)$$

and when $P > 0$,

$$T = \frac{\sum_{i=1}^N p_i \cdot t_i}{P}. \quad (4)$$

Sometimes a block may have one or more edges that loop back to itself. Often what is of interest is the total probability and expected cost to transition from one block to a different block. In these cases, it is useful to merge the effects of a loop-back edge into the edges that lead away from a block by appropriately adjusting the probabilities and expected costs of those edges that lead away. Assume a block has an edge that leads away with probability p and expected cost t and a loop-back edge with probability p_{loop} and expected cost t_{loop} . Then, the total probability P and expected cost T that a transition is made away from the block via the particular away edge being considered, accounting for an arbitrary number of possible loops before making

that transition is given by

$$P = \begin{cases} \frac{p}{1-p_{loop}} & \text{for } p_{loop} < 1, \\ 0 & \text{for } p_{loop} = 1 \end{cases} \quad (5)$$

and when $P > 0$,

$$T = t + \frac{p_{loop}}{1-p_{loop}} \cdot t_{loop}. \quad (6)$$

9.3 Base States

For the base computation and recovery states, p_0 is the probability that the application executes for a certain period of time without encountering a failure, and t_0 is the length of this period. Further, p_k for $k \in 1, 2, \dots, L$, is the probability that a level k failure occurs before any other failure during this period, and t_k is the expected run time before encountering such a failure. Given $f_k(t)$ and $F_k(T)$, we can directly compute the elements of \vec{p} and \vec{t} for the base states. First, we show the general solution. Assume there are L different classes, or *levels*, of failures, numbered $1, 2, \dots, L$. Assume a failure in one class is independent of failures in all other classes, and assume failures within a class are independent of one another. Since failures at different levels are assumed to be independent, the probability that there are no failures at any level during the time interval $t = 0$ to $t = T$ is given by

$$p_0(T) = (1 - F_1(T)) \cdot (1 - F_2(T)) \cdots (1 - F_L(T)),$$

and the expected run time given that no failures occur during that interval is simply

$$t_0(T) = T.$$

We next compute the probability that a level i failure will occur *before* a failure occurs at any other level during the time interval $t = 0$ to $t = T$. Consider a very small interval of time from t to $t + \Delta t$, where $0 \leq t < T$. The probability that a level i failure will occur during this interval is approximately $f_i(t) \cdot \Delta t$. The probability that a failure at another level has not occurred up until time t is

$$\begin{aligned} & (1 - F_1(t)) \cdot (1 - F_2(t)) \cdots \\ & \cdots (1 - F_{i-1}(t)) \cdot (1 - F_{i+1}(t)) \cdots \\ & \cdots (1 - F_L(t)). \end{aligned}$$

Furthermore, for each level $k \in 1, 2, \dots, L$, the probability that a level k failure will occur *before* a failure

occurs at any other level during the time interval $t = 0$ to $t = T$ is given by

$$p_k(T) = \int_0^T (1 - F_1(t)) \cdot (1 - F_2(t)) \cdots \\ \cdots (1 - F_{k-1}(t)) \cdot f_k(t) \cdot (1 - F_{k+1}(t)) \cdots \\ \cdots (1 - F_L(t)) dt.$$

The integrand above expresses the probability that a level k failure will occur during some infinitesimally small interval of width dt starting at time t ,

$$f_k(t) dt,$$

multiplied by the probability that a failure at another level has not already occurred by time t ,

$$(1 - F_1(t)) \cdot (1 - F_2(t)) \cdots \\ \cdots (1 - F_{k-1}(t)) \cdot (1 - F_{k+1}(t)) \cdots \\ \cdots (1 - F_L(t)).$$

The integral then sums the probabilities of each of these small intervals for all values of t between 0 and T to derive the total probability that a level k failure will occur before a failure occurs at any other level during the full time interval from $t = 0$ to $t = T$.

Similarly, when $p_k(T) > 0$, the expected run time given that a level k failure occurs before a failure occurs at any other level during the time interval $t = 0$ to $t = T$ is given by

$$a = \int_0^T t \cdot (1 - F_1(t)) \cdot (1 - F_2(t)) \cdots \\ (1 - F_{k-1}(t)) \cdot f_k(t) \cdot (1 - F_{k+1}(t)) \cdots (1 - F_L(t)) dt \\ t_k(T) = \frac{a}{p_k(T)}.$$

The expressions for $p_0(T)$ and $t_0(T)$ evaluate to $p_0(T) = e^{-\lambda T}$ and $t_0(T) = T$, and for $k \in 1, 2, \dots, L$, $p_k(T)$ and $t_k(T)$ evaluate to

$$p_k(T) = \frac{\lambda_k}{\lambda} (1 - e^{-\lambda T}), \quad t_k(T) = \frac{1 - (\lambda T + 1) \cdot e^{-\lambda T}}{\lambda \cdot (1 - e^{-\lambda T})},$$

where $\lambda = \lambda_1 + \lambda_2 + \dots + \lambda_L$.

9.3.1 Computation state: a $Y(k = 1, c)$ block

A $Y(k = 1, c)$ block is a base computation state in which the application executes for an interval of length t and then writes a checkpoint at level c , which requires a time of c_c . We denote the probability and

expected run time vectors for this state as \vec{p}_Y and \vec{t}_Y , respectively. Using the formulas from Section 9.3 we find that $p_{Y0} = p_0(t + c_c)$ and $t_{Y0} = t_0(t + c_c)$, and for $i \in 1, 2, \dots, L$, $p_{Yi} = p_i(t + c_c)$, and when $p_{Yi} > 0$, $t_{Yi} = t_i(t + c_c)$.

9.3.2 Recovery state at level k

While in a recovery base state at level k , the system is recovering from a failure using a checkpoint saved at level k , which requires a time of r_k . We denote the probability and expected run time vectors for this state as \vec{p}_R and \vec{t}_R . Using the formulas from Section 9.3, substituting i for k as a subscript label, and setting $T = r_k$, we find that $p_{R0} = p_0(r_k)$ and $t_{R0} = t_0(r_k)$.

If a failure occurs while in recovery, the recovery must be restarted. We specify which recovery state the system transitions to when discussing the X block. For now, we simply fill in \vec{p}_R and \vec{t}_R with values based on the recovery state required to handle the failure that occurs while in recovery, ignoring the level of the recovery state the system is currently in when that failure occurs. Thus, for each failure level $i \in 1, 2, \dots, L$, the probability that a level i failure occurs before a failure at any other level during the interval from $t = 0$ to $t = r_k$ is

$$p_{R0} = p_0(r_k) = e^{-\lambda \cdot r_k}, \quad t_{R0} = t_0(r_k) = r_k,$$

and for $i \in 1, 2, \dots, L$, $p_{Ri} = p_i(r_k)$, and when $p_{Ri} > 0$, $t_{Ri} = t_i(r_k)$.

9.4 The $X(k, c)$ block

An $X(k, c)$ block internally consists of a $Y(k, c)$ block and a recovery state at level k , which we denote by R_k . We refer to these three components as simply X , Y , and R . A diagram of an $X(k, c)$ block is shown in Figure 13(a). Here we compute the probability and expected run time vectors for X , which we denote as \vec{p}_X and \vec{t}_X , assuming that we are given the vectors for Y , as \vec{p}_Y and \vec{t}_Y , and R , as \vec{p}_R and \vec{t}_R . To simplify the final expressions, we merge groups of related transitions into single transitions. We show the merged transitions in Figure 13(b) (See Section 9.2 for details on merging transitions.)

Y transitions to the recovery state R for any failure scenario that requires a recovery level at k or less. We merge each of these transitions from Y to R into a single transition having probability of P_{YR} and an expected run time of T_{YR} , as shown in the top portion of Figure 14. Using formulas 3 and 4 from Section 9.2,

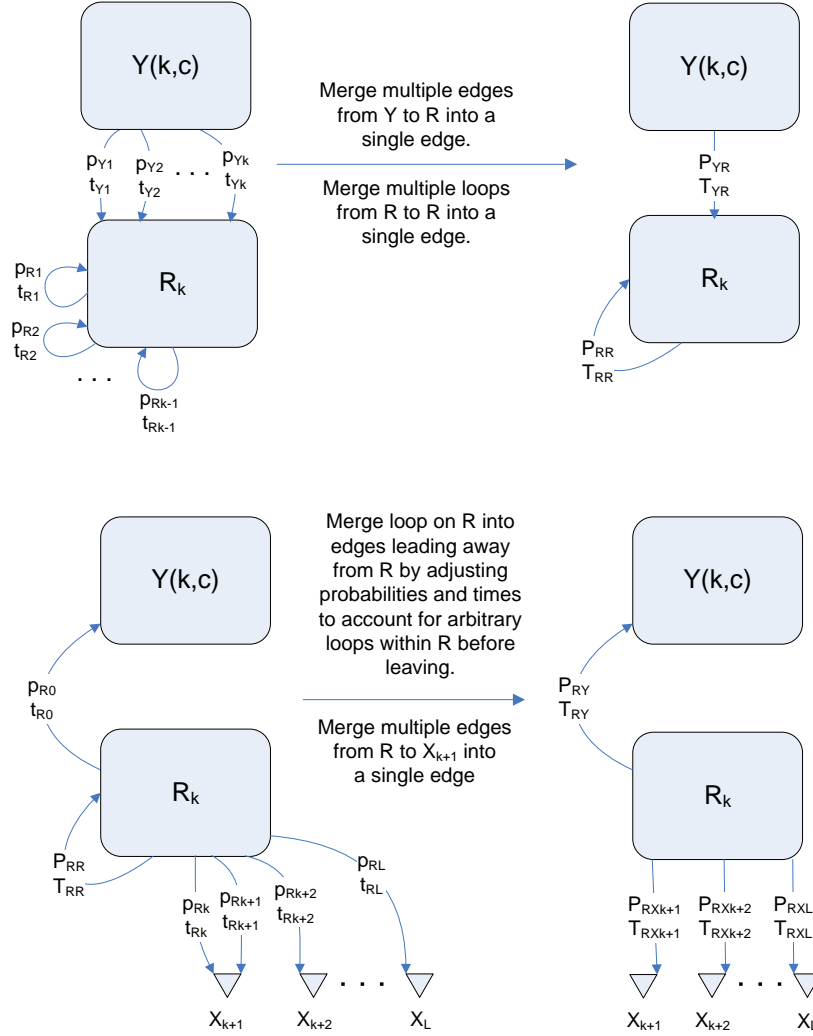


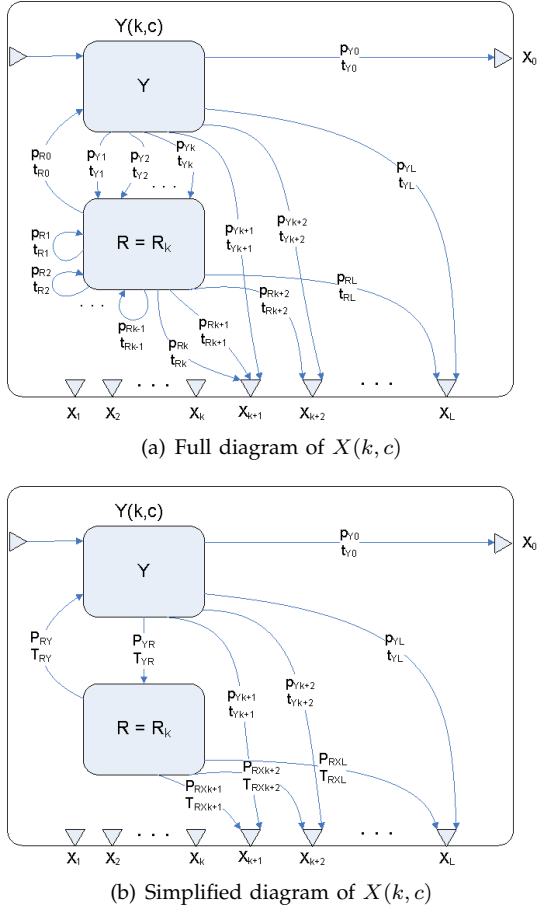
Fig. 14. Merging edges in $X(k, c)$

we get

$$P_{YR} = \sum_{i=1}^k p_{Yi} \quad T_{YR} = \frac{\sum_{i=1}^k p_{Yi} \cdot t_{Yi}}{P_{YR}}.$$

The R state has zero or more edges that loop back to itself, because a recovery state at the maximum level $k = L$ is restarted upon the occurrence of a failure at any level, and a recovery state at a level $k < L$ is restarted upon the occurrence of a failure at any level less than k . When $k < L$, a failure occurring at

level k or higher requires a transition from R_k to a higher level recovery state. In particular, while in R_k , we assume a failure at level k requires the system to transition to a recovery state of at least level $k + 1$. Again using formulas from Section 9.2, we merge each of these possible transitions from R to R into a single transition having probability P_{RR} and an expected run

Fig. 13. The $X(k, c)$ State

time T_{RR} as shown in the top portion of Figure 14:

$$P_{RR} = \begin{cases} 0 & \text{for } k = 1, \\ \sum_{i=1}^{k-1} p_{Ri} & \text{for } 1 < k < L, \\ \sum_{i=1}^k p_{Ri} & \text{for } k = L \end{cases}$$

and, when $P_{RR} > 0$,

$$T_{RR} = \begin{cases} \frac{\sum_{i=1}^{k-1} p_{Ri} \cdot t_{Ri}}{P_{RR}} & \text{for } 1 < k < L, \\ \frac{\sum_{i=1}^k p_{Ri} \cdot t_{Ri}}{P_{RR}} & \text{for } k = L. \end{cases}$$

The last substitution we make is to collapse the single loop transition from R to R into the transitions leading away from R . Upon entering R , a transition away from R eventually happens, provided $P_{RR} < 1$. However, one or more loops back to R may occur before transitioning away. It is possible to adjust the probabilities and expected run times of the transitions leading away from R to account for an arbitrary number of loops back to R before making the transition away. Derivation of the formulas to make these adjustments is

provided in Section 9.2. We collapse this loop transition into the transitions leading away from R by defining new transitions out of R having probabilities P_{RY} and P_{RX_i} with expected run times T_{RY} and T_{RX_i} for $i \in k, k+1, \dots, L$ as shown in the bottom portion of Figure 14. First, consider the transition from R to Y . Using formulas from Section 9.2, we collapse the loop from R to Y define its probability to be P_{RY} and expected run time to be T_{RY} as shown in Figure 13(b). First, we collapse the loop into the transition from R to Y by adjusting the probability of this transition from p_{R0} to P_{RY} with to P_{RY} . We account for the effects of the loop on this transition by adjusting its probability and expected run time by collapsing the loop into this transition we get

$$P_{RY} = \begin{cases} \frac{p_{R0}}{1-P_{RR}} & \text{for } P_{RR} < 1, \\ 0 & \text{for } P_{RR} = 1, \end{cases}$$

and, when $P_{RY} > 0$,

$$T_{RY} = \begin{cases} t_{R0} & \text{for } P_{RR} = 0, \\ t_{R0} + \frac{P_{RR}}{1-P_{RR}} \cdot T_{RR} & \text{for } 0 < P_{RR} < 1. \end{cases}$$

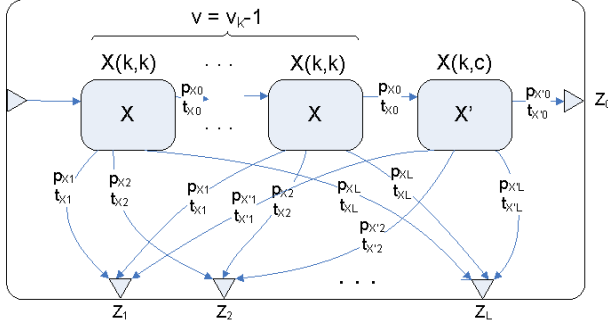
The transitions from R to external recovery states have probability P_{RX_i} and expected run time T_{RX_i} , as shown in Figure 13(b). Recall under our model that a recovery state at the maximum level $k = L$ is restarted upon the occurrence of a failure at any level. Such a state only contains loop-back transitions. In this case, there is zero probability that R transitions to an external recovery state, so when $k = L$, we have for each $i \in 1, 2, \dots, L$ $P_{RX_i} = 0$. While in a recovery state at level k , where $k < L$, the system transitions to a recovery state at level $k+1$ if either a level k failure or a level $k+1$ failure occurs. Otherwise, for the occurrence of a failure at level i , where $i > k+1$, a transition is made to a recovery state at level i . Thus, when $k < L$, we have for each $i \in 1, 2, \dots, L$

$$P_{RX_i} = \begin{cases} 0 & \text{for } 1 \leq i \leq k \text{ or } P_{RR} = 1, \\ \frac{p_{Rk} + p_{R(k+1)}}{1-P_{RR}} & \text{for } i = k+1 \text{ and } P_{RR} < 1, \\ \frac{p_{Ri}}{1-P_{RR}} & \text{for } i > k+1 \text{ and } P_{RR} < 1, \end{cases}$$

and, when $P_{RX_i} > 0$,

$$T_{RX_i} = \begin{cases} \frac{p_{Rk} \cdot t_{Rk} + p_{R(k+1)} \cdot t_{R(k+1)}}{p_{Rk} + p_{R(k+1)}} & \text{for } i = k+1, \\ t_{Ri} + \frac{P_{RR}}{1-P_{RR}} \cdot T_{RR} & \text{for } i > k+1. \end{cases}$$

After making all of these transformations, the simplified $X(k, c)$ block is shown in Figure 13(b). Finally, we derive the total probabilities and expected run times to

Fig. 15. The $Z(k, c)$ State

transition out of the X block as

$$p_{X0} = \begin{cases} \frac{p_{Y0}}{1 - P_{YR} \cdot P_{RY}} & \text{for } P_{YR} \cdot P_{RY} < 1, \\ 0 & \text{for } P_{YR} \cdot P_{RY} = 1, \end{cases}$$

and, when $p_{X0} > 0$,

$$t_{X0} = t_{Y0} + \frac{P_{YR} \cdot P_{RY}}{1 - P_{YR} \cdot P_{RY}} \cdot (T_{YR} + T_{RY}).$$

Also, for each level $i \in 1, 2, \dots, L$

$$p_{Xi} = \begin{cases} 0 & \text{for } 1 \leq i \leq k \text{ or } P_{YR} \cdot P_{RR} = 1, \\ \frac{p_{Yi} + P_{YR} \cdot P_{RXi}}{1 - P_{YR} \cdot P_{RY}} & \text{for } i > k \text{ and } P_{YR} \cdot P_{RR} < 1, \end{cases}$$

and, when $p_{Xi} > 0$,

$$t_{Xi} = \frac{p_{Yi} \cdot t_{Yi} + P_{YR} \cdot P_{RXi} \cdot (T_{YR} + T_{RXi})}{p_{Yi} + P_{YR} \cdot P_{RXi}} + \frac{P_{YR} \cdot P_{RY}}{1 - P_{YR} \cdot P_{RY}} \cdot (T_{YR} + T_{RY}).$$

9.5 The $Z(k, c)$ block

A $Z(k, c)$ block only exists when $v_k > 0$, and consists of a chain of $X(k, k)$ blocks of length $v_k - 1$ followed by a single $X(k, c)$ block, shown in Figure 15. We refer to these blocks as simply Z , X , and X' , where $Z = Z(k, c)$, $X = X(k, k)$, and $X' = X(k, c)$. Also, we define v such that $v = v_k - 1$. Here we compute the probability and expected run time vectors for Z , which we denote as \vec{p}_Z and \vec{t}_Z , assuming that we are given the vectors for X , as \vec{p}_X and \vec{t}_X , and X' , as $\vec{p}_{X'}$ and $\vec{t}_{X'}$.

The probability of successfully transitioning from the Z block to the first computation state of the next block is the probability that v consecutive successful transitions out of X blocks are followed by one successful transition out the X' block, $p_{Z0} = (p_{X0})^v \cdot p_{X'0}$. When $p_{Z0} > 0$, the expected time to make this transition is $t_{Z0} = v \cdot t_{X0} + t_{X'0}$.

For each failure level $i \in 1, 2, \dots, L$, there are multiple paths through Z which require a transition

to a recovery state at level i . The first X block may transition to a recovery state at level i . Or, that block may transition successfully to the next X block, which in turn may transition to a recovery state at level i . Or, the first two X blocks may transition successfully on to the the third X block, which may transition to a recovery state at level i , and so on up to and including the final X' block. The total probability to leave Z for a recovery state at level i is the sum of the probabilities corresponding to each of the possible paths from the substates of Z ,

$$\begin{aligned} p_{Zi} &= p_{Xi} + p_{X0} \cdot p_{Xi} + (p_{X0})^2 \cdot p_{Xi} \\ &\quad + \dots + (p_{X0})^{v-1} \cdot p_{Xi} + (p_{X0})^v \cdot p_{X'i} \\ p_{Zi} &= (1 + p_{X0} + (p_{X0})^2 \\ &\quad + \dots + (p_{X0})^{v-1}) \cdot p_{Xi} + (p_{X0})^v \cdot p_{X'i}. \end{aligned}$$

which can be simplified to

$$p_{Zi} = \begin{cases} \frac{1 - (p_{X0})^v}{1 - p_{X0}} \cdot p_{Xi} + (p_{X0})^v \cdot p_{X'i} & \text{for } p_{X0} < 1, \\ p_{X'i} & \text{for } p_{X0} = 1. \end{cases}$$

When $p_{Zi} > 0$, the time to transition from Z to a recovery state at level i is

$$t_{Zi} = \frac{A(p_{X0}, t_{X0}, p_{Xi}, t_{Xi}, p_{X'i}, t_{X'i})}{p_{Zi}}$$

where

$$\begin{aligned} A(p_0, t_0, p, t, p', t') &= B(p_0, t_0, p, t) + (p_0)^v \cdot p' \cdot (v \cdot t_0 + t'), \\ B(p_0, t_0, p, t) &= (1 + (p_0)^1 + (p_0)^2 + \dots + (p_0)^{v-1}) \\ &\quad \cdot p \cdot t + (1 \cdot (p_0)^1 + 2 \cdot (p_0)^2 + \dots + \\ &\quad (v-1) \cdot (p_0)^{v-1}) \cdot p \cdot t_0. \end{aligned}$$

9.6 The $Y(k, c)$ block

A $Y(k, c)$ block is built using three different constructions depending on the values of k and (when $k > 1$) v_{k-1} . If $k = 1$, then $Y(k, c) = Y(k = 1, c)$, which is a base computation state. The probability and expected run time vectors for this state can be directly computed as described in Section 9.3.1.

If $k > 1$ and $v_{k-1} = 0$, then $Y(k, c)$ consists of a single $Y(k - 1, c)$ block. We refer to these blocks as Y and Y' , where $Y = Y(k, c)$ and $Y' = Y(k - 1, c)$. Here we compute the probabilities and expected run times to transition out of Y given the probabilities and expected run times to transition out of Y' as That is, we compute \vec{p}_Y and \vec{t}_Y given $\vec{p}_{Y'}$ and $\vec{t}_{Y'}$. Because Y consists solely of Y' , the probability and expected run time vectors to transition out of Y , which we denote as \vec{p}_Y and \vec{t}_Y , are trivially computed given the vectors

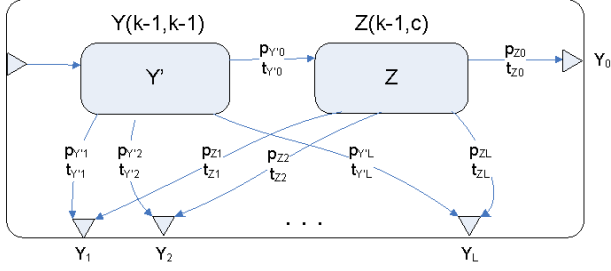


Fig. 16. The $Y(k, c)$ State for $k > 1$ and $v_{k-1} > 0$

for Y' , as $\vec{p}_{Y'}$ and $\vec{t}_{Y'}$: $p_{Y_0} = p_{Y'0}$ and $t_{Y_0} = t_{Y'0}$, and, for each level $i \in 1, 2, \dots, L$, $p_{Y_i} = p_{Y'i}$ and $t_{Y_i} = t_{Y'i}$.

If $k > 1$ and $v_{k-1} > 0$, then $Y(k, c)$ consists of a starting $Y(k-1, k-1)$ block followed by a $Z(k-1, c)$ block, as shown in Figure 16. We refer to these states as Y , Y' , and Z , where $Y = Y(k, c)$, $Y' = Y(k-1, k-1)$, and $Z = Z(k-1, c)$. The probability and expected run time vectors to transition out of Y , which we denote as \vec{p}_Y and \vec{t}_Y , assuming that we are given the vectors for Y' , as $\vec{p}_{Y'}$ and $\vec{t}_{Y'}$, and Z , as \vec{p}_Z and \vec{t}_Z . The probability that a successful transition out of Y occurs is the probability that both Y' and Z transition successfully, $p_{Y_0} = p_{Y'0} \cdot p_{Z0}$ and the expected time for this transition is $t_{Y_0} = t_{Y'0} + t_{Z0}$.

For each failure level $i \in 1, 2, \dots, L$, there are two possible paths through Y that can cause a transition to a recovery state at level i . The Y' block may transition immediately to a recovery state at level i , or the Y' block may transition successfully to Z , which in turn may transition to a recovery state at level i . The total probability to leave Y for a recovery state at level i is the sum of the probabilities corresponding to each path $p_{Y_i} = p_{Y'i} + p_{Y'0} \cdot p_{Zi}$. When $p_{Y_i} > 0$, the expected run time of this transition is

$$t_{Y_i} = \frac{p_{Y'i} \cdot t_{Y'i} + p_{Y'0} \cdot p_{Zi} \cdot (t_{Y'0} + t_{Zi})}{p_{Y_i}}.$$

10 SCAVENGE MODEL

Here, we extend our model to explore the scenario where level L checkpoints are written only upon job termination instead of periodically throughout the run. We refer to the process of pushing a checkpoint set to the parallel file system upon allocation termination as *scavenging a checkpoint*. Scavenging only occurs on application failure. Level $k < L$ checkpoints will be taken at regular intervals during the application run and cached on the compute nodes.

Our approach to computing efficiency with our model is based on computing the amount of work done

in each state. On failure, if it is possible to scavenge a checkpoint from a state, we export the expected amount of work done before the failure occurs. If it is not possible to scavenge a checkpoint from a state, we export the expected amount of work done on failure transitions only. A transition out of a state on success means that no failures occurred and the amount of work accomplished is the work accumulated in any substates of the current state.

If a level k failure occurs and the current state contains a level $i \geq k$ checkpoint, we can scavenge a checkpoint; we call this *transitioning on scavenge*. In this case, the amount of work completed is the amount of work accumulated in substates before the failure. A transition on failure occurs if the current state does not contain a level $i \geq k$ checkpoint or a level $i \geq k$ failure occurs during the scavenge process. On failure transitions at level $i > k$, the amount of work completed in the current state is zero. In the extreme case, if no completed checkpoints can be scavenged and the application must roll back to the last level L checkpoint, the total expected amount of work completed is zero.

Changes to our model are: Recovery states within $X(k, c)$ are now scavenge states. Transitions from each compute state at failure levels $i \leq k$ are now called scavenge transitions. Instead of maintaining the probability, expected total time, and expected amount of work at each failure level for scavenge transitions, we compute the overall expected values.

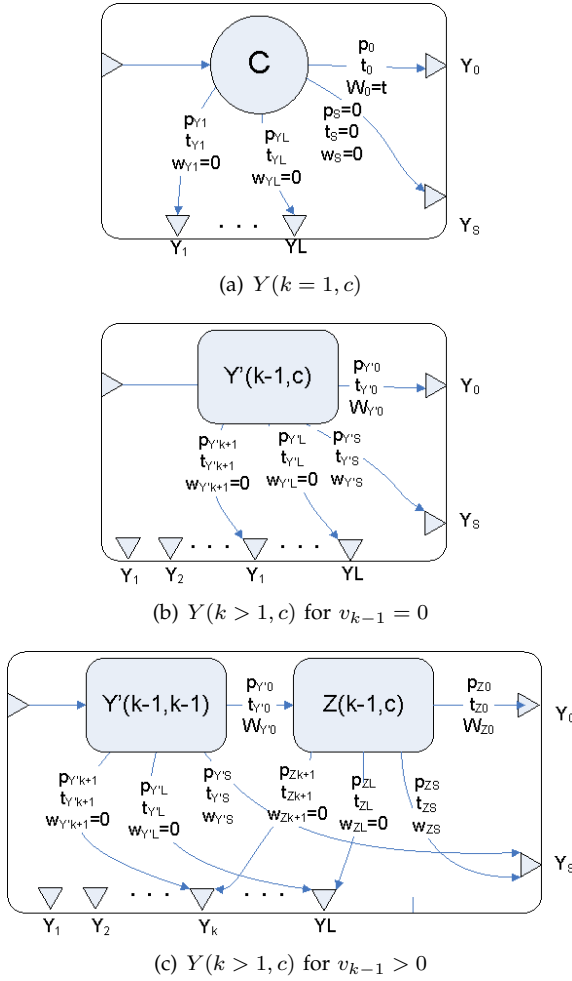
10.1 Y States

10.1.1 Computation state: a $Y(k=1, c)$ block

In the simplest case, we have a $Y(k=1, c)$ block, which only contains a single compute interval of duration t followed by a level c checkpoint, shown in Figure 17(a). Upon transition on success from this state, the expected amount of work completed will be t and the expected amount of time in Y is t_0 , as derived in Section 9.3.1. Upon transition from Y on failure for any level $i \in 1, 2, \dots, L$, $w_{Y_i} = 0$ because no work interval and checkpoint were completed successfully. The probabilities and times for transitioning out of Y are the same as were derived in Section 9.3.1. There is no possibility of transitioning out of Y on scavenge, because there is no scavenge state contained within a $Y(k=1, c)$ state. Therefore, on transition out of Y on scavenge, $p_S = 0$, $t_S = 0$, and $w_S = 0$.

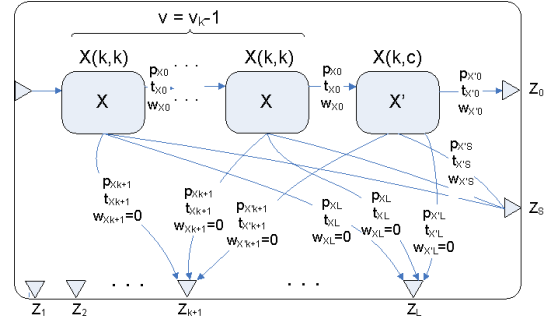
10.1.2 The $Y(k > 1, c)$ block

In the case where $v_{k-1} = 0$, Y consists of a single enclosed $Y(k-1, c)$ state, Y' as shown in Figure 17(b).

Fig. 17. The $Y(k, c)$ States

The probabilities and times for transitioning out on success or failure are computed in Section 9.6. The work completed for a transition out of Y on success or failure is $w_{Y0} = w_{Y'0}$ and $w_{Yi} = w_{Y'i} = 0$. For a scavenge transition, we have $p_{YS} = p_{Y'S}$, $t_{YS} = t_{Y'S}$, and $w_{YS} = w_{Y'S}$.

In the case where $v_{k-1} > 0$, Y consists of a $Y(k-1, k-1)$ state, Y' , and a $Z(k-1, c)$ state, Z (Figure 17(c)). Here, on transition out of Y on success, p_{Y0} and t_{Y0} are the same as derived in Section 9.6 and w_{Y0} is the sum of the work done in Y' and Z , $w_{Y0} = w_{Y'0} + w_{Z0}$. On a transition out of Y on failure at levels $i \in 1 \dots L$, then again p_{Yi} and t_{Yi} are the same as derived in Section

Fig. 18. $Z(k, c)$ State

9.6 and $w_{Yi} = 0$. On a transition out of Y on scavenge,

$$p_{YS} = p_{Y'S} + p_{Y'0} \cdot p_{ZS}$$

$$t_{YS} = \frac{p_{Y'S} \cdot t_{Y'S} + p_{Y'0} \cdot p_{ZS}(t_{Y'0} + t_{ZS})}{p_{YS}}$$

$$w_{YS} = \frac{p_{Y'S} \cdot w_{Y'S} + p_{Y'0} \cdot p_{ZS}(w_{Y'0} + w_{ZS})}{p_{YS}}.$$

10.2 Z State

A $Z(k, c)$ state, Z consists of $v = v_k - 1$ consecutive $X(k, k)$ states X , followed by a single $X(k, c)$ state X' (Figure 18). The probabilities and times for transitioning out of Z on success are the same as derived in Section 9.5, and the expected amount of work that is performed in Z on successful transition out of the state is simply $w_{Z0} = w_{X0} \cdot v + w_{X'0}$. On transition out of Z on failure for $i \in 1 \dots L$, p_{Zi} and t_{Zi} remain the same as in Section 9.5 and the amount of work completed is $w_{Xi} = 0$.

The probability of a transition out of Z on scavenge is the probability that any of the subset X states transition on scavenge. Using the result from Section 9.5 for p_{Zi} , we have

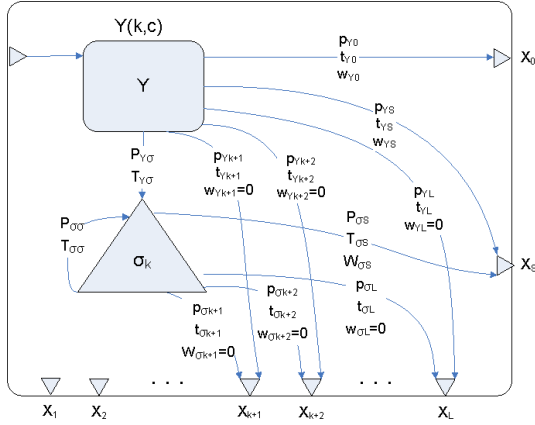
$$p_{ZS} = \begin{cases} \frac{1-(p_{X0})^v}{1-p_{X0}} \cdot p_{XS} + (p_{X0})^v \cdot p_{X'S} & \text{for } p_{X0} < 1 \\ p_{X'S} & \text{for } p_{X0} = 1. \end{cases}$$

Similarly, the expected time before leaving Z on scavenge is

$$t_{ZS} = \frac{A(p_{X0}, t_{X0}, p_{XS}, t_{XS}, p_{X'S}, t_{X'S})}{p_{ZS}},$$

where A is defined in Section 9.5. And finally, the amount of work performed in Z before transition out on scavenge is

$$w_{ZS} = \frac{A(p_{X0}, w_{X0}, p_{XS}, w_{XS}, p_{X'S}, w_{X'S})}{p_{ZS}}.$$

Fig. 19. $X(k, c)$ State

10.3 X State

An $X(k, c)$ state, X consists of a $Y(k, c)$ state, Y , and a scavenge state σ_k . A failure at level $i \leq k$ in Y causes a transition from Y to σ_k . For any failures at levels $i < k$ that occur in σ_k , the scavenge operation restarts. Failures at levels $i \geq k$ in Y or $i \geq k$ in σ_k , cause a transition out of X on failure.

The probability and time for transitioning out of X on success are $p_{X0} = p_{Y0}$ and $t_{X0} = t_{Y0}$. The expected amount of work that is performed in X on successful transition out of the state is the amount of work that was performed in Y , $w_{X0} = w_{Y0}$.

A failure at level $i \leq k$ in Y causes a transition to the scavenge state σ_k . The amount of work done on transition from Y to σ is $w_{Y\sigma} = 0$. The probability and time for transitioning from Y to σ_k on failure at level $i \leq k$ are

$$P_{Y\sigma} = \sum_{i=1}^k p_{Yi} \quad T_{Y\sigma} = \frac{\sum_{i=1}^k p_{Yi} \cdot t_{Yi}}{P_{Y\sigma}}.$$

As stated previously, we estimate the time for a successful scavenge operation as the time for a level L checkpoint, c_L . From Section 9.3, the probabilities and times before transitioning out of σ_k on success or failure at level i are

$$p_{\sigma 0} = p_0(c_L) \quad p_{\sigma i} = p_i(c_L) \quad t_{\sigma 0} = t_0(c_L) \quad t_{\sigma i} = t_i(c_L).$$

Failures at levels $i < k$ cause a restart of the scavenge operation. The probability and time before a restart of the scavenge operation are $P_{\sigma\sigma}$ and $T_{\sigma\sigma}$. No work is done in σ_k , so $w_{\sigma i} = 0$. The expected probability and time before a level $i < k$ failure and restart of the

scavenge operation are

$$P_{\sigma\sigma} = \sum_{i=1}^{k-1} p_{\sigma i} \quad T_{\sigma\sigma} = \frac{\sum_{i=1}^{k-1} p_{\sigma i} \cdot t_{\sigma i}}{P_{\sigma\sigma}}.$$

Given the derivations for merging loop-back edges in Section 9.2, The probability and time for leaving σ_k on success or failure for $i \in 0, 1, \dots, L$ are

$$P_{\sigma i} = \frac{p_{\sigma 0}}{1 - P_{\sigma\sigma}} \quad T_{\sigma i} = t_{\sigma 0} + \frac{P_{\sigma\sigma}}{1 - P_{\sigma\sigma}} \cdot T_{\sigma\sigma}.$$

The amount of work accomplished is $W_{\sigma i} = w_{\sigma 0}$. However, if $P_{\sigma\sigma} = 1$ then, $P_{\sigma i}$, $T_{\sigma i}$, and $W_{\sigma i}$ are zero. The probability, time, and amount of work accomplished on leaving σ_k for a scavenge state at level i or higher for $i \geq k$ are

$$P_{\sigma i} = \frac{p_{\sigma i}}{1 - P_{\sigma\sigma}} \quad T_{\sigma i} = t_{\sigma i} + \frac{P_{\sigma\sigma}}{1 - P_{\sigma\sigma}} \cdot T_{\sigma\sigma} \\ W_{\sigma i} = 0.$$

When transitioning out of X on failure at level $i = k$, the probability, time, and work only depend on the scavenge state. However, when transitioning out of X on failure at level $i > k$, the values depend on the Y and scavenge states,

$$P_{Xi} = \begin{cases} P_{Y\sigma} \cdot P_{\sigma i} & \text{for } i = k \\ p_{Yi} + P_{Y\sigma} \cdot P_{\sigma i} & \text{for } i > k. \end{cases}$$

$$T_{Xi} = \begin{cases} \frac{P_{Y\sigma} \cdot P_{\sigma i} (T_{Y\sigma} + T_{\sigma i})}{p_{Yi} \cdot t_{Yi} + P_{Y\sigma} \cdot P_{\sigma i} (T_{Y\sigma} + T_{\sigma i})} & \text{for } i = k \\ \frac{P_{Xi}}{P_{Xi}} & \text{for } i > k \end{cases}$$

$$W_{Xi} = 0$$

The probabilities, times, and amount of work done for transitions out of X on successful scavenge are

$$P_{XS} = p_{YS} + P_{Y\sigma} \cdot P_{\sigma S} \\ T_{XS} = \frac{p_{YS} \cdot t_{YS} + P_{Y\sigma} \cdot P_{\sigma S} (T_{Y\sigma} + T_{\sigma S})}{P_{XS}} \\ W_{XS} = \frac{p_{YS} \cdot w_{YS} + P_{Y\sigma} \cdot P_{\sigma S} (w_{\sigma S})}{P_{XS}}.$$